

Detecting People in Multi-view Environment using Height Map Reconstruction in Real-Time

Ákos Kiss

Distributed Events Analysis Laboratory, Computer and Automation Research Institute (MTA SZTAKI), H-1111, Kende u. 13-17, Budapest, Hungary

Dept. of Control Engineering and Information Technology, Budapest University of Technology and Economics, H-1117 Budapest, Magyar tudósok körútja 2., Budapest, Hungary

Tamás Szirányi

Distributed Events Analysis Laboratory, MTA SZTAKI

Abstract

In this article we address the problem of multiview detection: we detect pedestrians by their feet in scene space. Detection is based on searching cone intersections. The altitude of detection is also retrieved, which eliminates the need of planar ground - which is a common restriction in the related literature. We found positions can be computed accurately, and despite a large number of false positives, height map of the scene can be reconstructed with small error. Precision of the detector can be increased knowing height map, so that results of our method are comparable to SOA methods in case of planar ground - using somewhat different test methodology -, but adding the ability to handle arbitrary ground. Our algorithm is capable of real-time operation, based on two optimizations: decreasing number of cones, and approximating intersection bodies. Cones are back-projections of ellipses in images covering feet, this is ensured by filtering and clustering foreground pixels. Moreover, most demanding steps are parallelizable, and distributable due to lack of data dependencies.

Keywords: multi-view detection, 3D position, projection, real-time processing

1. INTRODUCTION

Single camera detecting and tracking relies on some kind of descriptors, like color, shape, and texture - a survey about state of the art in monocular pedestrian detection can be found in [3]. However, extraordinary and occluding

Email addresses: kiss.akos@sztaki.mta.hu (Ákos Kiss), sziranyi.tamas@sztaki.mta.hu (Tamás Szirányi)

15 objects are hard to detect. Multiple cameras are often used to overcome these limitations.

16 In these cases, consistency of object pixels among views will signal objects in 3D space. There are many methods
17 for finding object pixels. Using stereo cameras, the disparity map can highlight foreground, or using wide-baseline
18 stereo imaging, image-wise foreground detection is carried out.

19 The resulting set of object pixels can be further segmented, and consistency check might be reduced to likely
20 correspondent segments. This might be done using color descriptors [14, 15], however color calibration [10] is
21 necessary due to different sensor properties or illumination conditions.

22 Foreground masks can be projected to a ground plane, and overlapping pixels mark consistent regions [9]. Robust-
23 ness can be improved using multiple parallel planes [12]. Furthermore, looking for certain patterns in the projection
24 plane increases reliability [16].

25 In many works authors assume known homography between views and ground plane to carry out projection. In
26 [7] homography parameters are estimated from co-motion statistics from multimodal input videos, eliminating the
27 need of human supervision.

28 Projection of whole foreground masks is computationally expensive, but filtering pixels can reduce complexity. In
29 some works, points associated with feet are searched, reducing foreground masks from arbitrary blobs to points and
30 lines [13, 9].

31 Another gain of filtering foreground points is that - depending on the geometry - feet are less occluded than whole
32 bodies, eliminating a great source of errors. Reducing occlusion is especially important in dense crowds, for example
33 using top-view cameras [5].

34 **2. OVERVIEW**

35 Our goal was to design an algorithm for detecting people in a multiview environment with possibly many views
36 where the geometry of the ground is arbitrary, which problem is not addressed in the related literature. Moreover we
37 aimed at reaching real-time operation to make it available for surveillance systems.

38 We utilized the fact that foreground pixels correspond to lines in 3D scene space, and these lines intersect at
39 surface points of the object. Finding such intersections however would require too many computations, as number of

40 line pairs n_p is proportional to

$$n_p = (s_x s_y r_{fg})^2 \binom{n}{2} \quad (1)$$

41 where s_x, s_y is the resolution of the camera image, r_{fg} is the ratio of foreground pixels in image and n is the number
42 of views. We found that increasing number of views is not the bottleneck, but the great number of foreground pixels.
43 Number of these pixels increases with resolution, and errors in foreground mask. Therefore we chose to

- 44 1. filter the foreground mask to find pixels relevant to detecting position of people, and
- 45 2. collect spatially coherent pixels and form one primitive from them instead of handling many lines.

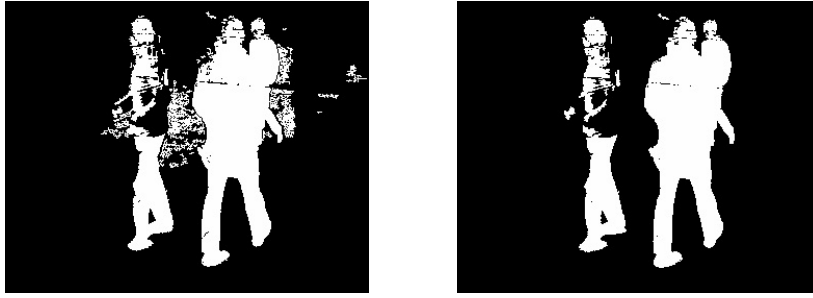
46 We filtered pixels possibly corresponding to feet, which we call candidate pixels. We covered these pixels with
47 ellipses, these can be back-projected to cones in scene space. Intersecting a pair of cones replaces intersecting many
48 line pairs.

49 Our approach has several advantages in means of both precision and speed:

- 50 • for determining cone parameters, undistortion may be carried out with few computations leading to accurate
51 parameters,
- 52 • number of cones is proportional to number of objects regardless of image resolution, this makes our cone
53 matching approach scalable compared to pixel matching,
- 54 • ground doesn't have to be flat (unlike using homographies in [16, 12, 11, 9, 2]),
- 55 • we can compute synergy map without height range presumption like in [16, 12, 11].

56 On the downside:

- 57 • our algorithm may fail on incorrect foreground mask. As we will show in Section 3.2, if foreground mask is
58 incorrect, the extracted ellipses won't cover feet precisely, or won't even intersect any foot.



(a) Extracting foreground mask without (left) and with (right) white balance compensation.



(b) Extracting candidate pixels of feet from foreground mask.

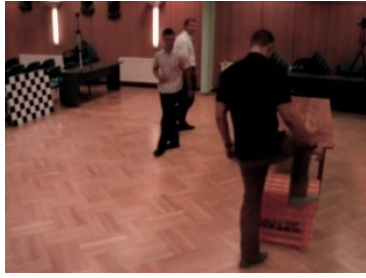


(c) Candidate pixels are clustered, clusters are modeled with ellipses (ellipses are visualized with lozenges).

Figure 1: Steps of extracting feet

- 59 • precise calibration is required for reliable estimation of cone parameters (both intrinsic and extrinsic parameters
60 of the camera). A small error in cone or camera parameters can slightly change the direction of the cone, which
61 results in a great error far from the camera.

62 We show applied preprocessing steps to form ellipses from foreground mask in section 3. Section 4 introduces
63 steps of forming cones in scene space and matching these cones. Section 5 describes feet detection from cone matches



(a) Shadows may corrupt foreground detection even in carefully chosen color space. Foreground mask (right) for an input frame (left) is shown. A number of false positives are introduced due to corrupted foreground mask.



(b) Candidate pixels can appear on arms or on foreground artifacts, and also cones corresponding to different feet can intersect.



(c) Corresponding images showing time skew.

Figure 2: Sources of errors

64 as well as details of height map reconstruction. We show results and comparison to state of the art methods in Section
 65 6 and conclude our work in Section 7.

66 3. PREPROCESSING

67 We found that camera’s auto white balance function leads to incorrect foreground mask. This function aims at
 68 adopting to changes in lightning, but it is also affected by large objects entering or exiting the image, thereby changing
 69 white balance even if lightning conditions remain the same.

70 To neglect these effects, we extended the background model described in [1] with white balance adaptation. By
 71 computing and neglecting consistent changes in background pixel values, we managed to compensate changes in
 72 white balance. This resulted in a more reliable foreground detection (see Fig 1a)

73 Foreground detection was corrupted by shadows and reflections on the ground. We tried different color spaces to
 74 overcome this issue, and chose to use XYZ color space. It lead to best overall performance of our algorithm, however
 75 in certain situations foreground mask was still corrupted by shadows, as we can see in Fig 2a.

76 We filtered out small areas from the foreground mask to suppress noise. In our tests, we eliminated areas covering
77 less than 7 pixels, as we found foreground blobs due to noise did not exceed this size. For other videos, this threshold
78 might be adjusted according to video quality. Note that this step is not necessary, it only speeds up further processing.

79 *3.1. Filtering foreground pixels*

80 We assumed that cameras are in upright position so that pixels of feet are bottom pixels of vertical lines in the
81 foreground mask. We consider these pixels as candidates for feet, we call them candidate pixels. A sample of extracted
82 candidate pixels can be seen in Fig 1b.

83 Note that using calibration information, images are easily rotated or search in a different upright direction may be
84 carried out.

85 We clustered the candidate pixel set by connecting pixels closer than 3px. This is necessary, as candidate pixels
86 are often not connected because of noise and steep edges. Note that these effects don't scale with resolution, so this
87 threshold can be used for different image sizes, however, adjustment might be required for different image quality.

88 Clustering results in a set of pixel sets. We model each pixel set with an ellipse according to the spatial mean
89 and variance of the pixels. This makes our algorithm robust to image resolution, as increasing resolution results in
90 more candidate pixels, but the same number of clusters - thus reducing the number of 3D primitives, which drastically
91 speeds up pairwise matching. Only the steps of preprocessing is affected by the resolution.

92 *3.2. Forming Ellipses*

93 A cluster is considered to be a sampling of an elliptical area. For an ellipse with major and minor radius a and b
94 parallel to axis, we know:

$$\int x^2 = \frac{a^3 b \pi}{4} \quad (2)$$

$$\int y^2 = \frac{a b^3 \pi}{4} \quad (3)$$

$$\int xy = 0 \quad (4)$$

95 Rotating with α we get:

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix} \quad (5)$$

$$\int x'^2 = c^2 \int x^2 + s^2 \int y^2 \quad (6)$$

$$\int y'^2 = s^2 \int x^2 + c^2 \int y^2 \quad (7)$$

$$\int x'y' = sc \left(\int x^2 - \int y^2 \right) \quad (8)$$

$$\text{tg}2\alpha = \frac{2 \int x'y'}{\int y'^2 - \int x'^2} \quad (9)$$

$$\int x'^2 + \int y'^2 = \int x^2 + \int y^2 = \frac{ab\pi}{4} (a^2 + b^2) \quad (10)$$

$$\frac{\int x'y'}{sc} = \int x^2 - \int y^2 = \frac{ab\pi}{4} (a^2 - b^2) \quad (11)$$

$$\int 1 = ab\pi \quad (12)$$

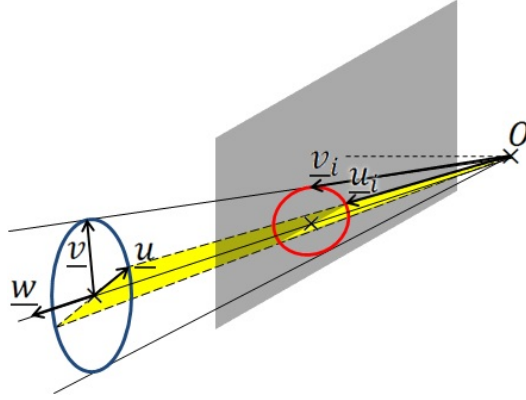
96 Where $s = \sin\alpha$, $c = \cos\alpha$. Now we can compute major and minor radius and rotation. We reject too small and
 97 upright ellipses. Minimal size can be adjusted for different setups and image resolution, in our experiments ellipses
 98 covering less than 7 candidate pixels were rejected, this ensured feet far from the camera are still detected. We defined
 99 maximal angle as 45° , foot direction does not depend on setup or resolution, so this threshold should work in any
 100 setups. Sample results can be seen in Fig. 1c.

101 4. DETECTING CORRESPONDENCES

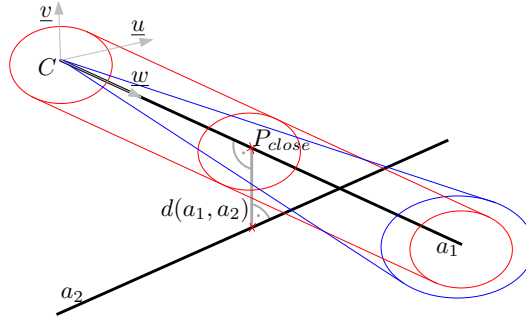
102 Cones corresponding to a foot will intersect close to the location of foot. However, cones corresponding to
 103 different feet may accidentally intersect leading to false detections. Moreover, candidate pixels can appear on arms or
 104 on foreground artifacts, introducing false cones. Fig 2b shows examples of such false matches.

105 4.1. Forming Cones

106 We chose to describe a cone with a vertex C and three orthogonal vectors \underline{u} , \underline{v} and \underline{w} , where \underline{w} is the unit length



(a) Model of a cone. Vertex C is in the projection center O , \underline{w} points toward center of ellipse in image plane. Bevel angles and axes directions (\underline{u} , \underline{v}) are computed from axes in image plane (\underline{u}_i , \underline{v}_i).



(b) Cylinder for cone matching.

Figure 3: Steps of extracting feet

107 direction vector of axis and \underline{u} , \underline{v} are direction vectors of major and minor axes (Fig. 3a). Bevel angles are determined
 108 by the length of \underline{u} , \underline{v} vectors. Cone consists of points \underline{p} where

$$\begin{aligned} ((\underline{p}-C)\underline{u})^2 + ((\underline{p}-C)\underline{v})^2 &\leq ((\underline{p}-C)\underline{w})^2, \\ (\underline{p}-C)\underline{w} &\geq 0 \end{aligned} \tag{13}$$

109 From extrinsic parameters of camera we know 3D position of points of image plane as well as the optical center
 110 O . \underline{w} points from O to the center of the ellipse, and $\overrightarrow{Ou_i}$, $\overrightarrow{Ov_i}$ are in symmetry planes, where u_i and v_i are extremal
 111 points of ellipse in image plane (Fig 3a). We determined major and minor axes directions by vector products to ensure
 112 orthogonality:

$$\underline{v} = \underline{u}_i \times \underline{w} \quad (14)$$

$$\underline{u} = \underline{w} \times \underline{v} \quad (15)$$

113 Finally, \underline{u} and \underline{v} are scaled according to major and minor bevel angles so that (13) stands.

114 4.2. Cone Matching

115 Intersection of cones is a complex body, especially when bevel-angles differ for major and minor axes - apart
 116 from some extraordinary situations. Determining this body would be difficult, but our experiments showed that we
 117 don't need the exact intersection, an approximate solution is sufficient. We applied simple heuristics for finding cone
 118 intersections - which we call matches. We simplified computation in two ways:

- 119 1. generally a foot covers a small area of the image, so cones will have small bevel-angles. Consequently in a
 120 relatively small space segment, they can be replaced by elliptical cylinders. We assume that intersection will
 121 take place near the point where the two cone axes are closest, so we compute major and minor radii at this offset
 122 $\overrightarrow{CP_{close}}$ in Fig. 3b).
- 123 2. Computing exact intersection of cylinders is still a hard problem, and the intersection body is still complex.
 124 Therefore instead of computing intersection body, we tried to find an optimal point \underline{p} in space, for which
 125 distance from axes of the elliptical cylinder is minimal - considering different radii.

126 After computing $|\overrightarrow{CP_{close}}|$ distance, we scale \underline{u} and \underline{v} vectors to match major and minor cylinder radii:

$$\begin{aligned} \underline{u}' &= \underline{u} |\overrightarrow{CP_{close}}| \\ \underline{v}' &= \underline{v} |\overrightarrow{CP_{close}}| \end{aligned} \quad (16)$$

127 Now, equation of the elliptical cylinder becomes

$$((\underline{p} - C)\underline{u}')^2 + ((\underline{p} - C)\underline{v}')^2 \leq 1 \quad (17)$$

128 Distance from axis - left side of (17) - is a linear function of \underline{p} , enabling us to write a linear equation system
 129 expressing \underline{p} is on both axes (index refers to cylinders):

$$\begin{aligned} ((\underline{p}-C_1)\underline{u}'_1)^2 + ((\underline{p}-C_1)\underline{v}'_1)^2 &= 0 \\ ((\underline{p}-C_2)\underline{u}'_2)^2 + ((\underline{p}-C_2)\underline{v}'_2)^2 &= 0 \end{aligned} \tag{18}$$

equivalent to

$$\begin{bmatrix} \underline{u}'_1{}^T \\ \underline{v}'_1{}^T \\ \underline{u}'_2{}^T \\ \underline{v}'_2{}^T \end{bmatrix} \underline{p} = \begin{pmatrix} \underline{u}'_1 C_1 \\ \underline{v}'_1 C_1 \\ \underline{u}'_2 C_2 \\ \underline{v}'_2 C_2 \end{pmatrix} \tag{19}$$

$$A\underline{p} = \underline{b} \tag{20}$$

130 Of course, axes will practically never intersect, only approximate solution is possible. Solving subject to least
 131 square error is straightforward and practical in our case. It is straightforward, because it is a distance from axes
 132 considering different major and minor radii, and it is practical, because this optimization is easy to carry out - for
 133 example using pseudo inverse.

134 If \underline{p} is outside any cylinder, we conclude cones are not intersecting, otherwise, they intersect and \underline{p} is the position
 135 of the intersection - which we call match. We also assign a weight to matches computed from the error of the
 136 approximate solution.

137 5. DETECTING FEET

138 5.1. Detection

139 One match itself can be a false detection. However, for two matches, foot has to be visible in at least three views,
 140 which is not common in dense crowd. Therefore we applied a threshold to the weight of matches to decide if foot is
 141 detected.

142 We merged matches close in space to avoid multiple detections. We also assigned weight to such sets, computed
143 from the weights of the matches in a way, that the possibility of detection highly increases with the number of matches.

144 To ease merging close matches, we merged matches over a grid on horizontal plane. As feet don't appear above
145 each other, in a cell of the grid we don't mix matches at different altitudes. We chose to collect matches in blocks,
146 which consist of overlapping 2×2 cells, then compute the weight of match sets in every block. Using overlapping
147 blocks is necessary to achieve smooth weight map over the grid.

148 Every local maximum in this 2D weight map with a weight above a given threshold is considered a detected foot.
149 This threshold is a parameter of our algorithm, higher threshold leads to increased precision along with decreased
150 recall. 3D position of foot is computed as the barycenter of matches in the block, where z coordinate is height
151 (parallel to upright direction).

152 5.2. Height Map

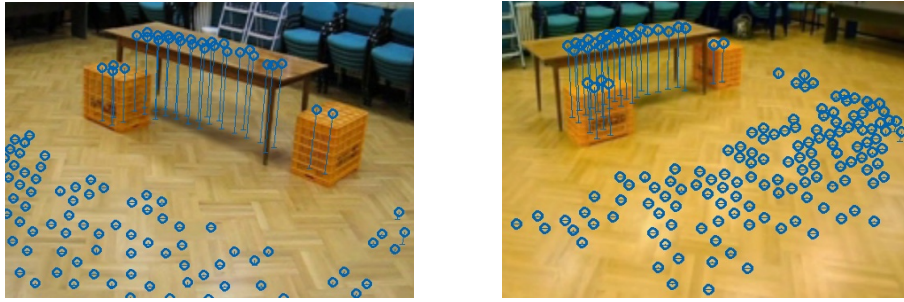
153 As we mentioned before, and showed in Fig 2b, false detection may occur due to incorrect foreground detection
154 or random intersections of non-corresponding cones. We found that in a dense crowd many false detections appear.
155 However the height of these detections is quite random, and they usually occur far from the ground.

156 Consequently these false positives appear as a noise. To suppress this noise, we aggregated all detections from a
157 test set and applied statistical filtering with normal distribution presumption. As a result, we got positions on ground
158 with high confidence. These dense points will be height points of ground.

159 The phrase height map is somewhat misleading, because at the border of surfaces at different altitude, it is possible
160 to have height points above each other, however in our tests this never occurred. Thereby we kept the height map
161 expression.

162 After extracting height map, we can ignore detections far from any height point, leading to a much more reliable
163 method. We speeded up height point searching by accumulating detections in a 3D grid, which makes periodic rerun
164 possible. This enables us to refine height map as more and more detections occur.

165 A sample height map of our non planar test case is shown in Fig 4a. There are certain areas where few detections
166 took place, this led to incomplete height map.



(a) Generated height map points for our test case.



(b) Corresponding frames from different cameras with the detected feet marked

Figure 4: Sample results of detection and determining height map.

167 6. EXPERIMENTS

168 We tested our algorithm on a commonly used test sequence for multiview detection, the *EPFL terrace* sequence
 169 [4] (sample results can be seen in Fig 4b), and our own test videos (*SZTAKI* sequence).

170 We made test videos of a scene where three planar surfaces are present in different heights. This was required to
 171 demonstrate capabilities of our method as available multiview test sets present planar ground. We recorded videos

172 from 4 views with different kind of consumer digital cameras with video capture function. We found that our algorithm
173 performed well despite the different camera parameters, distortion and image quality.

174 Our approach, detecting foot has advantages and disadvantages as well. Feet are always near ground, which makes
175 height map detection possible and helps rejecting false positives. Also, for certain camera setups, feet are less likely
176 to be occluded compared to whole bodies.

177 On the other hand in case of detecting foot positions, precise time synchronization is necessary, as a small time
178 skew can corrupt detection. In Fig. 2c, we can see the moving leg is bended in one view while straight in the other.
179 Hardly noticeable in the small image, however the distance in space is greater than the size of the foot itself.

180 6.1. Height Map Reconstruction

181 We reconstructed height map for both *EPFL* and *SZTAKI* sequences. We found that height map reconstruction
182 was carried out with high precision in both cases.

183 6.1.1. EPFL sequence

184 In case of *EPFL* sequence, we measured that the reconstructed ground is a flat plane with small error: $\mu = -1.6cm$
185 (mean), $\sigma = 1.7cm$ (std. deviation). Fig 5a (left histogram) shows the histogram of height values.

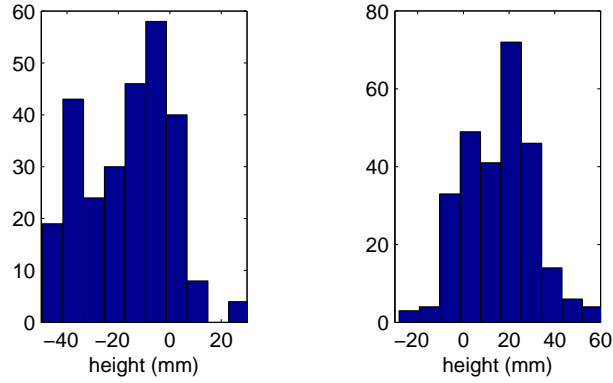
186 We measured distribution different from normal distribution - which we expected assuming noise-like error. A
187 possible source of systematic error is incorrect plane normal, so we computed optimal normal vector by fitting a plane
188 on height points. The optimal normal vector was very close to upright direction specified by camera calibration. It
189 diverged by only 0.36° , resulting in slightly smaller error: $\mu = 1.6cm$, $\sigma = 1.5cm$. However, using this normal led to
190 a significantly better distribution (Fig 5a, right histogram).

191 6.1.2. SZTAKI sequence

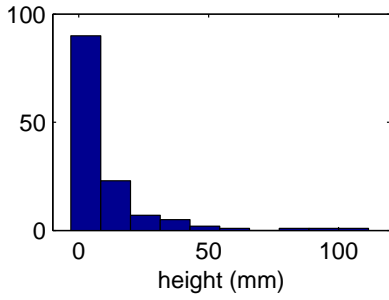
192 In case of our sequence, we found that all three surfaces were found, and their height was determined with high
193 accuracy. Measurements are summarized in Table 1a. Fig 5b shows the histogram of heights for points on floor.

194 There are 3 outlier points (above 6cm high) for the floor. These appear on the edge of the area of interest, as can
195 be seen in Fig 4a. We found this is due to sagging people - which lead to false candidate pixels on image borders.

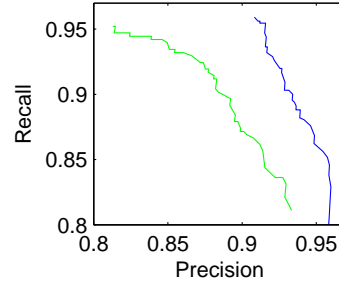
196 Without these outlier points, mean and σ for floor becomes 0.9cm and 1.1cm respectively.



(a) Histogram of height of ground points for *EPFL* dataset in reference coordinate system - according to calibration data (left), and in direction of computed normal vector (right).



(b) Histogram of height of floor points for *SZTAKI* dataset.



(c) ROC curve measured in function of detection threshold on *EPFL* (blue), and *SZTAKI* (green) sequences.

Figure 5: Evaluation of detection and height point calculation.

197 6.2. Detecting People

198 We tested our method using manually created ground truth information of feet positions. We implemented an
 199 application to create ground truth, in which we have to mark a foot in at least 2 views and its position is determined
 200 by solving a linear equation similar to (19). If we mark foot in more views, number of rows of A and length of \underline{b}
 201 increases.

202 We collected feet positions frame by frame, person by person. 1 or 2 legs can be specified for a person, because
 203 sometimes only one foot is visible from more views.

204 For the evaluation, we defined a region of interest (ROI), a rectangle on the ground plane and matched detections
 205 to ground truth inside this area. In case of *EPFL* dataset, this rectangle is specified along with the dataset, for our test
 206 case this area was chosen so that every part is visible from at least 3 views.

Table 1: Numerical evaluation of two main aspects: accuracy and running time.

(a) Statistical information on surfaces found in scene.

surface	height	μ	σ	# of points	min	max
floor	0cm	1.1cm	1.7cm	131	-0.3cm	11.2cm
box	50cm	50.2cm	0.6cm	6	49.1cm	50.8cm
table top	73cm	74.4cm	1.2cm	23	70.3cm	75.9cm

(b) Average processing time of each step for all 4 views (single threaded implementation, 2.4GHz Core 2 Quad CPU). Foreground detection, filtering and cone forming can be done parallel (per camera), cone matching and detection requires all cone information to be present.

test set	foreground detection	candidate filtering	cone forming	Σ (per camera)	cone matching	detection	Σ (per system)
<i>EPFL</i>	51.2ms	3.37ms	63us	13.7ms	879us	28us	907us
<i>SZTAKI</i>	32.5ms	2.88ms	1.99ms	9.35ms	600us	18us	618us

207 We accepted a detection as true positive, if it was not further than 25cm from a ground truth foot position. This is
 208 approximately the length of a foot. Afterwards we eliminated people and detections outside the ROI, and computed
 209 false negatives and false positives from left people and detections accordingly.

210 As mentioned in Section 5.1, a detection threshold is applied to the weight map. Trade off between precision and
 211 recall is balanced with this threshold. To evaluate our method, we ran the test with different threshold values, resulting
 212 in precision-recall curves. Results are shown in Fig 5c.

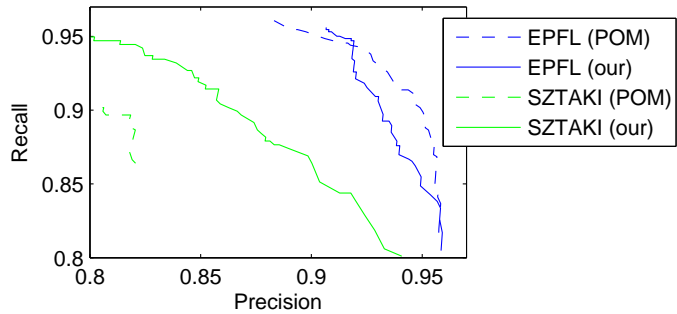
213 6.2.1. Comparison to SoA Method

214 We compared our method to two SoA methods referred to as POM[6] and 3DMPP[16]. Comparison is not
 215 straightforward, as these methods detect bodies projected to a ground plane, in contrast we detect feet positions in
 216 space. Consequently, a different evaluation criteria had to be used: we consider someone detected if either of his legs
 217 were detected.

218 Fig 6a shows results found in [16] for POM and 3DMPP - evaluated on EPFL and PETS datasets - compared to

	Precision	Recall	Time (without foreground detection)
POM	87.20%	95.56%	~0.5s
3DMPP	97.5%	95.5%	~1s
our	90.66%	95.61%	~0.01s

(a) Precision/recall values measured on EPFL dataset. Results for POM and 3DMPP are copied from [16].



(b) ROC curve measured in function of detection threshold on *EPFL* (blue), and *SZTAKI* (green) sequences.

Figure 6: Comparing our results to a SoA methods.

219 our measurements - evaluated only on EPFL dataset. We used similar amount of test data (179 frames - every 25th
 220 -, 661 objects), and the same ROI. Results show that our method is comparable to these SoA methods in case of flat
 221 ground.

222 We also compared results on our dataset to show effect of non-planar ground. For this we chose POM as our
 223 framework was easy to generate input for it's available implementation.

224 For our dataset, foreground detection was challenging due to reflective ground and poor video quality. As a result,
 225 both methods performed worse in this case, however our algorithm performed noticeably better (see Fig 6b).

226 6.3. Running Time

227 We measured average running time of processing steps, results are shown in Table 1b. Every step highly depends
 228 on input data, thus ratio of foreground mask processing times differ from ratio of pixels. In these cases, the video
 229 resolution was 360×288 for *EPFL* and 320×240 for *SZTAKI* sequences

230 Results show that our method is capable of real-time processing of 4 views with 360×288 resolution even with
 231 single threaded implementation. However, parallel processing is also possible with some restrictions.

232 Several steps of processing perform on single images, therefore can be parallelized. Clearly, foreground detection,
 233 filtering, ellipse and cone forming is done imagewise. These steps can be done even on different computers, or on
 234 smart cameras, as no data exchange is required between views at all.

235 Cone matching and detection requires all cone information to be present. Due to short processing time, these steps
 236 should run sequentially on a dedicated resource. In a camera network, where not all cameras have overlapping field

237 of view, camera sets can be processed individually.

238 In case of above mentioned SOA methods, mean processing time of foreground mask is in the order of 0.5s for
239 POM, and 1s for 3DMPP, using unoptimized code on similar processor (data provided by author of [16]). However in
240 our case, processing foreground took less than 10ms.

241 Moreover, many works - including 3DMPP - operate projecting parts or whole foreground masks to planes, which
242 makes distributed computing impossible due to data dependencies.

243 7. CONCLUSION

244 We proposed a multiview-detection algorithm that retracts 3D position of people using multiple calibrated and
245 synchronized views. In our case, unlike other algorithms, non-planar ground can be present. This is done by modeling
246 possible positions of feet with 3D primitives, cones in scene space and searching for intersections of these cones.

247 For good precision, height map of ground should be known. Our method can compute height map on the fly,
248 reaching high precision after a startup time.

249 After height map detection we measured precision and recall values comparable to SOA methods on a commonly
250 used data set. Our algorithm worked well also on our test videos we made to demonstrate capabilities of handling
251 non-planar ground.

252 In the future we plan to examine tracking people by their leaning leg positions[8].

253 8. ACKNOWLEDGE

254 The authors would like to thank François Fleuret for making implementation of POM available, and Ákos Utasi
255 for sharing his experiments of evaluating POM and 3DMPP algorithms - his code is also available for download.

256 This work has been supported by the Hungarian Scientific Research Fund grant OTKA #106374.

257 9. REFERENCES

258 [1] Benedek, C. and Szirányi, T. (2008). Bayesian foreground and shadow detection in uncertain frame rate surveillance videos. *IEEE Image*
259 *Processing*, 17(4):608–621.

- 260 [2] Berclaz, J., Fleuret, F., and Fua, P. (2006). Robust people tracking with global trajectory optimization. In *IEEE CVPR*, pages 744–750.
- 261 [3] Dollar, P., Wojek, C., Schiele, B., and Perona, P. (2012). Pedestrian detection: An evaluation of the state of the art. *IEEE Trans. Pattern Anal.*
262 *Mach. Intell.*, 34(4):743–761.
- 263 [4] EPFL (2011). Multi-camera pedestrian videos. <http://cvlab.epfl.ch/data/pom/>.
- 264 [5] Eshel, R. and Moses, Y. (2010). Tracking in a dense crowd using multiple cameras. *International Journal of Computer Vision*, 88:129–143.
- 265 [6] Fleuret, F., Berclaz, J., Lengagne, R., and Fua, P. (2008). Multicamera people tracking with a probabilistic occupancy map. *IEEE Trans.*
266 *Pattern Anal. Mach. Intell.*, 30(2):267–282.
- 267 [7] Havasi, L. and Szlavik, Z. (2011). A method for object localization in a multiview multimodal camera system. In *CVPRW*, pages 96–103.
- 268 [8] Havasi, L., Szlavik, Z., and Szirányi, T. (2007). Detection of gait characteristics for scene registration in video surveillance system. *IEEE*
269 *Image Processing*, 16(2):503–510.
- 270 [9] Iwase, S. and Saito, H. (2004). Parallel tracking of all soccer players by integrating detected positions in multiple view images. In *IEEE ICPR*,
271 pages 751–754.
- 272 [10] Jeong, K. and Jaynes, C. (2008). Object matching in disjoint cameras using a color transfer approach. *Machine Vision and Applications*,
273 19:443–455.
- 274 [11] Khan, S. and Shah, M. (2006). A multiview approach to tracking people in crowded scenes using a planar homography constraint. In *ECCV*
275 *2006*, Lecture Notes in Computer Science, pages 133–146.
- 276 [12] Khan, S. and Shah, M. (2009). Tracking multiple occluding people by localizing on multiple scene planes. *PAMI*, 31(3):505–519.
- 277 [13] Kim, K. and Davis, L. S. (2006). Multi-camera tracking and segmentation of occluded people on ground plane using search-guided particle
278 filtering. In *ECCV*, pages 98–109.
- 279 [14] Mittal, A. and Davis, L. (2001). Unified multi-camera detection and tracking using region-matching. In *IEEE Multi-Object Tracking*, pages
280 3–10.
- 281 [15] Mittal, A. and Davis, L. (2002). M2tracker: A multi-view approach to segmenting and tracking people in a cluttered scene using region-based
282 stereo. In *ECCV 2002*, pages 18–33.
- 283 [16] Utasi, Á. and Benedek, C. (2011). A 3-D marked point process model for multi-view people detection. In *IEEE CVPR*, pages 3385–3392.